

Conditionals in Homomorphic Encryption and Machine Learning Applications

Chialva, Diego; Dooms, Ann

Published in:
ArXiv.org

Publication date:
2018

[Link to publication](#)

Citation for published version (APA):

Chialva, D., & Dooms, A. (2018). Conditionals in Homomorphic Encryption and Machine Learning Applications. *ArXiv.org, 2018*.

Copyright

No part of this publication may be reproduced or transmitted in any form, without the prior written permission of the author(s) or other rights holders to whom publication rights have been transferred, unless permitted by a license attached to the publication (a Creative Commons license or other), or unless exceptions to copyright law apply.

Take down policy

If you believe that this document infringes your copyright or other rights, please contact openaccess@vub.be, with details of the nature of the infringement. We will investigate the claim and if justified, we will take the appropriate steps.

Conditionals in Homomorphic Encryption, and Machine Learning Applications

Diego Chialva and Ann Dooms

Abstract—Homomorphic encryption aims at allowing computations on encrypted data without decryption other than that of the final result. This could provide an elegant solution to the issue of privacy preservation in data-based applications, such as those using machine learning, but several open issues hamper this plan. In this work we assess the possibility for homomorphic encryption to fully implement its program without relying on other techniques, such as multiparty computation (SMPC), which may be impossible in many use cases (for instance due to the high level of communication required). We proceed in two steps: i) on the basis of the structured program theorem [Bohm, Jacopini] we identify the relevant minimal set of operations homomorphic encryption must be able to perform to implement any algorithm; and ii) we analyse the possibility to solve -and propose an implementation for- the most fundamentally relevant issue as it emerges from our analysis, that is, the implementation of conditionals (requiring comparison and selection/jump operations). We show how this issue clashes with the fundamental requirements of homomorphic encryption and could represent a drawback for its use as a complete solution for privacy preservation in data-based applications, in particular machine learning. Our approach for comparisons is novel and entirely embedded in homomorphic encryption, while previous studies relied on other techniques, such as SMPC, demanding high level of communication among parties, and decryption of intermediate results from data-owners. A number of studies have indeed dealt with comparisons, but typically their algorithms rely on other techniques, such as secure multiparty computation, which required a) high level of communication among parties, and b) the data owner to decrypt intermediate results. Our protocol is also provably safe (sharing the same safety as the homomorphic encryption schemes), differently from other techniques such as Order-Preserving/Revealing-Encryption (OPE/ORE).

Index Terms—homomorphic encryption, machine learning

I. INTRODUCTION

Machine learning, data mining and predictive data analytics represent an ensemble of techniques and algorithms (which for simplicity we will in the following indicate simply as "machine learning") that allow systems to act and make predictions without being explicitly programmed in full detail to do so, but by leveraging their input data with inference techniques. They have nowadays an overwhelming number of practical applications providing us with an unprecedented level of comfort and services, from tailored suggestion systems, to "personalised medicine", and several other services.

D. Chialva is with the ERCEA (European Research Council Executive Agency) and A. Dooms is with the Department of Mathematics, Vrije Universiteit Brussel, Pleinlaan 2, 1050 Brussels, Belgium.

Disclaimer. The views expressed in this paper by Diego Chialva are the author's. They do not necessarily reflect the views or official positions of the European Commission, the European Research Council Executive Agency or the ERC Scientific Council.

However, these advantages typically come at the price of losing individual privacy, as personal or valuable information is used by the algorithms and the third parties operating them. This issue has spawned research activity at different levels. Very roughly speaking we can divide the developed privacy-preservation techniques in two classes: those that work by modifying the data themselves and those that modify the representation of the data, but not the actual data content.

Techniques of the first class act on the datasets holding the privacy-concerned data and can be divided in a few subclasses [2]. Common to all of them is the distinction between identifier, quasi-identifier and anonymous data. Such techniques require only comparatively minor (conceptual) changes to the application algorithms acting on the data, but they have significant drawbacks. Indeed, they impose a trade-off between the degree of preserved privacy and the usefulness of the data: a "privacy budget", which has been shown to be quite limited [28]. Moreover, such techniques appear to be beatable by the algorithms themselves and database crossing attacks (that is, the use by an attacker of other, public or stolen, databases to "complete" or infer the relevant distorted information in the database of interest) [28].

The other class of techniques has been proposed within cryptography. Among the different research lines we recall: secure multi-party computation, functional encryption, program obfuscation and homomorphic encryption, see for instance [7], [12], [29], [56]. These approaches differ in several aspects, including the set of functions that can be computed on the encrypted data and stage of development.

Homomorphic encryption aims at enabling the computation of arbitrary (in the case of *fully homomorphic encryption*) or classes (in the case of *partial homomorphic encryption*) of functions on encrypted data without having the need to decrypt them first and *limiting decryption to the very final result only*. This is in particular interesting for privacy preservation (including algorithm protection) in learning applications, and has been actively pursued in the latest years, e.g [6], [13], [14], [32], [35]. However, several open issues make homomorphic cryptosystems still unsuited for the vast majority of machine learning algorithms. Those that have been identified in the literature mainly are: memory footprint, computational complexity, limited representable data (only integers and finite precision floats) and a restricted set of operations (only polynomial operations, that is addition and multiplication).

Such problems can however be divided into two classes. The first class comprises issues such as the memory footprint and the computational complexity, which could be hoped to be trivially solved by technological (hardware) advancements,

similarly to what has happened in deep learning. On the other hand, the second class of the above-listed problems, like the limited types of representable data and the lack of more general operations, must find a solution at the theoretical and cryptography level. It is this second class of issues that we are interested in within this work.

We therefore individuate and analyse the minimum set of basic operations necessary to implement any algorithm, whatever its complexity, and assess the possibility/impossibility to implement them in homomorphic encryption from first principles.

We proceed on the basis of the well-known *structured program theorem* [9], which states that every computable function¹ can be implemented in a programming language that combines subprograms in only three specific ways:

- 1) Executing one subprogram, and then another subprogram (sequence);
- 2) Executing one of two subprograms according to the value of a Boolean variable (selection);
- 3) Executing a subprogram until a Boolean variable is true (iteration).

We observe that 2) and 3) require being able to perform *conditionals*, that is *comparison* operations to compare values and evaluate down to a single Boolean value, *and selection/jump* operations to pick up the correct branch of a program. Hence in order to assess the possibility for homomorphic encryption to accommodate all algorithms (in particular machine learning ones), comparisons and selections/jumps must be implementable².

In this work we address both the issue of comparison and of selection/jump operations. Concerning the former, several proposals have been made concerning comparison operations in an encrypted setting, but not yet totally within an homomorphic cryptosystem. There is even a claim that comparison would not be feasible in pure homomorphic encryption, see for example the comments in [4]. Typical, well-studied approaches have been

- 1) the so-called *Order-Preserving-Encryption* (OPE) and its variants such as *Modular-OPE* and *Order-Revealing-Encryption* (ORE), see for instance [3], [10], [11], [43], [46], [48], which do not belong to the homomorphic encryption class and, more seriously, have been proven to be not secure (for recent proofs, see for instance [8], [36]);
- 2) *secure multi-party computation* (SMPC) (for recent works see [22], [44]), which, although secure, require a high level of communication between parties (each single comparison in a machine learning training and prediction process must be performed by exchanging several messages), which may not be always possible;
- 3) combinations of homomorphic encryption with other cryptographic techniques such as SMPC to perform the

comparisons [16], [24], [41], [47], [55]. Again, these approaches do not manage to perform the comparisons exclusively on encrypted messages, as the data owner is required by the protocol to decrypt intermediate results, extract the significant bits for the comparison, re-encrypt and send the result back to the other party for the accomplishment of the algorithm. Such “decryption in the middle” hampers the purpose of homomorphic encryption (also, the need for a high level of communication between parties due to the use of SMPC may be impossible in a number of actual practical use cases).

In this paper we develop, through a new approach, a technique to achieve comparisons in homomorphic encryption (that is, with no need for communication between parties and acting exclusively on encrypted messages with no need for intermediate nor partial decryption).

When turning however to selection/jump operations, which are integral elements of conditionals, and thus of practically relevant algorithms, we will show that one hits a rather fundamental issue in homomorphic cryptosystems, namely the cryptographic requirement of *semantic security*. We will show how a more limited form of selection/jump operations can still be implemented, and we discuss the limitation the above-mentioned issue imposes on the implementation of full machine learning algorithms. In particular, this could represent a serious drawback in using homomorphic encryption for data analysis applications and for implementing algorithms in general, and could force to revisit (or abandon) that plan in its more ambitious formulation.

The article is organised as follows: we provide a brief introduction to homomorphic encryption in Section II, after which we study comparisons and selections in a homomorphic setting in Section IV. In Section V we present our methodology and test results. We finally discuss applications to machine learning in Sections VI, VI-A, where we highlight and provide precise and exhaustive examples of how the fundamental, general issues of homomorphic encryption that our analysis has revealed impact the program of implementing machine learning algorithms in such framework. In our conclusions we also briefly touch upon the consequences of our work in applications different than machine learning.

II. HOMOMORPHIC ENCRYPTION

A cryptosystem consists of three sets, a plaintext P , ciphertext C and key space K , together with a family of encryption functions $\text{Encr} : K \times P \rightarrow C$ and decryption functions $\text{Decr} : K \times C \rightarrow P$ such that for each $k \in K$, there exists a $k' \in K$ such that $\text{Decr}(k', \text{Encr}(k, p)) = p$ for all $p \in P$. Although in the literature Encr is called *encryption function*, it is not exactly a function in the strict mathematical sense for most of the encryption schemes, because an element of (pseudo)randomness is involved such that applying it more than one time to the same key and plaintext, one obtains different ciphertexts. Such probabilistic encryption schemes

¹Technically, representable as a flow chart, such as all machine learning algorithms.

²Proving the impossibility of implementing such operations, would entail the impossibility to implement complex algorithms under homomorphic encryption. Clearly, in the opposite case, where the fundamental analysis is positive, one still needs to assess the effectiveness of the implementation, which may still condemn homomorphic encryption to be impractical.

are favoured because they provide *semantic security*³, which is equivalent to *ciphertext indistinguishability*⁴ [33], [34]. This required randomness has a huge relevance in homomorphic encryption, as we will see.

Encryption schemes are further distinguished by the relation between the encryption and decryption key. If the decryption key can be easily computed from the encryption one (in the typical case they are in fact identical), one speaks of a *symmetric* cryptosystem, while if not, one speaks of an *asymmetric* cryptosystem. Typical asymmetric systems also distinguish between public (for encryption) and private (for decryption) keys k_p, k_s .

In modern cryptanalysis the adversaries are conceived as having finite computational resources and a cryptosystem is considered secure if its breaking is unfeasible with attack algorithms that are probabilistic in nature and running in polynomial time. The running of the cryptosystems functions and adversary algorithms are all measured as a function of the so-called *security parameter* λ , which measures the complexity of the computational problem.

A cryptosystem is *homomorphic* for an operation $*$ acting on P if there is a corresponding operation \circ acting on C with

$$\text{Decr}(k_s, \text{Encr}(k_p, m_1) \circ \text{Encr}(k_p, m_2)) = m_1 * m_2 \quad (1)$$

for $m_1, m_2 \in P$.

Note this is not in general a true group homomorphism, as

$$\text{Encr}(k_p, m_1) \circ \text{Encr}(k_p, m_2) \neq \text{Encr}(k_p, m_1 * m_2). \quad (2)$$

due to the (pseudo)randomness of the encryption scheme. However, while mathematically this lack of identity holds, there is a strong definition of homomorphic cryptosystems that reconciles with the group-homomorphism-like identity, in the statistical or computational senses, see [39].

Defining a *homomorphic encryption system* $(P, C, K, \text{Encr}, \text{Decr}, \text{Ev})$ then consists of specifying the evaluation function Ev that performs the homomorphically preserved operations O on (a number of) ciphertexts

$$\text{Ev} : C^n \times C \times K \rightarrow C : (\vec{c}, O, p_k) \rightarrow c' \quad (3)$$

where C is the family of circuits that the homomorphic cryptosystem can evaluate. An homomorphic cryptosystem is defined *correct* if it correctly decrypts ciphertexts both coming from a circuit evaluation (sometimes called “evaluated ciphertexts”), and from direct encryption of a plaintext (also dubbed “fresh ciphertexts”). Trivial homomorphic cryptosystems are excluded by requiring *strong homomorphicity* and *compactness* for the cryptosystems, for which we refer the reader to [39]. We will consider exclusively non-trivial homomorphic cryptosystems.

Homomorphic cryptosystem can be further distinguished in

³The fact that no polynomial time probabilistic algorithm can derive information about a plaintext m given its length, ciphertext and encryption algorithm, more than any other polynomial time probabilistic algorithm that has no access to the ciphertext.

⁴Given two plaintexts chosen by the adversary and the ciphertext of one of them chosen by us, the adversary cannot distinguish which of the plaintexts has been encrypted with a probability (significantly) larger than $1/2$, see [27].

- *partially homomorphic*: allow only one type of operation (addition or multiplication) for an *unlimited* number of times,
- *somewhat homomorphic*: allow addition and multiplication, but only for a *limited* number of times (the size of the ciphertext depends on the circuit depth),
- *levelled homomorphic*: allow addition and multiplication, but only for a *limited* number of times specified as an input parameter (here the size of the ciphertext does not depend on the maximal allowed circuit depth, but the size of the public key does),
- *fully homomorphic* allow addition and multiplication for an *unlimited* number of times, and thus arbitrary functions expressible as arithmetic circuits.

The randomness necessary for semantically secure schemes introduces a noise component that increases with each evaluation of an operation in the circuit. When the noise is above a certain limit, decryption is no longer correct. Fully homomorphic systems, as constructed first by Gentry, see [30], can cope with this issue thanks to a procedure, called *bootstrapping* that allows to extend specific somewhat homomorphic cryptosystems (called *bootstrappable*) to systems where unlimited number of operation evaluations are possible. Later realisations are, for example, [17], [21], [25], [31], [45], [53], [54]. For the purpose of this work, it is important to note that the increase in noise is different for addition and multiplication. Typically the one induced by multiplications is much larger.

Moreover, and quite relevantly, in practical applications the computational complexity (and hence slowness) of homomorphic operations has prompted to use *levelled* homomorphic systems. This clearly affects the algorithms that can be successfully implemented at the practical level.

III. APPROXIMATION BY POLYNOMIALS

By definition, somewhat, levelled and fully homomorphic encryption can only deal with polynomial operations, as they are the operations modelled by circuits of addition and multiplication. Therefore one takes the approach to approximate the functions one would like to evaluate homomorphically, with polynomials, e.g. see [40].

For the convenience of the reader, we report here some fundamental elements of the theory of polynomial approximations, which will be relevant in the following. It is well-known by the *Weierstrass Approximation Theorem* (see e.g. [51]) that any real-valued continuous function f on a closed interval $[a, b] \subset \mathbb{R}$ can be uniformly approximated by a polynomial, i.e. for every $\epsilon > 0$ there exists a polynomial p such that for all $x \in [a, b]$, $|f(x) - p(x)| < \epsilon$ or equivalently $\|f - p\|_\infty < \epsilon$, where

$$\|f\|_\infty = \sup\{|f(x)| : x \in [a, b]\} = \max\{|f(x)| : x \in [a, b]\},$$

the *supremum norm*. Moreover, there exists a unique polynomial of degree n that minimizes the supremum norm within the set of all polynomials of degree n , which is called the *polynomial of best approximation* or *minimax polynomial of degree n* . Being able to restrict the degree and to have control

over the maximum error makes this type of approximation very attractive. However, as the supremum norm is not induced by an inner product, the theory of orthogonal projections cannot be applied, but luckily there exist several (numerical) algorithms, such as the classical [49], that can determine the minimax polynomial.

IV. CONDITIONALS AND COMPARISONS IN HOMOMORPHIC CRYPTOSYSTEMS

Homomorphic encryption aims at computing any computable function on encrypted data *without recurring to intermediate, not even partial, decryption*, and it has been highly regarded as a possibility to make privacy-safe machine learning algorithms. As mentioned in the introduction, homomorphic encryption allows to compute general polynomial operations, but in order to apply homomorphic encryption at least conceptually to general algorithms one still needs to prove, on the basis of the structured program theorem, that it can provide for comparisons as well as selections, the two fundamental components of conditionals.

We will take a completely novel approach and show that some of the involved aspects pose rather crucial problems for the program of homomorphic encryption. Along the line of our proposed solution to implement comparisons and selections/jumps, we will also have to deal with other open issues in homomorphic encryption, such as the ability to perform divisions among ciphertexts.

A. Implementing Comparison Operations in Homomorphic Cryptosystems

We start by defining a comparison operation as a map

$$\text{Comp} : C \times C \rightarrow \mathcal{S} = \{0, \pm 1\} \quad (4)$$

where C is the ciphertext space. We tackle the problem by trying to find a representation of this map in terms of elements of the circuit family that the homomorphic cryptosystem can evaluate, i.e. polynomial operations, rather than trying to implement comparisons via additional basic/elementary features of our cryptosystem (as attempted in OPE/ORE, and so far inconclusive). However, Comp is **not** straightforwardly representable in terms of polynomials as it is discontinuous and typically implemented as a sign or equivalently using the *Heaviside (step) function*⁵,

$$H(x) = \begin{cases} 1 & \text{if } x > 1 \\ \frac{1}{2} & \text{if } x = 0 \\ 0 & \text{if } x < 0. \end{cases}$$

Note that $H(x) = \frac{1}{2}(1 + \text{sgn}(x))$.

Indeed, as H is discontinuous, the Weierstrass Approximation Theorem does not apply and insisting on such an approximation requires using many polynomials of high degree, while the approximations are still of bad quality because of Gibb's phenomenon. As high polynomial order implies a high number of consecutive multiplications in the homomorphic system,

this is problematic for the levelled or somewhat homomorphic schemes to which one is limited in practice as we have explained before.

We can however cope with these relevant issues and obtain a satisfactory definition and modelling.

Solution to the problem. We propose our solution, allowing: 1) to use only polynomial operations, 2) to compute comparisons in an efficient way in pure homomorphic encryption.

As we have remarked, Comp is typically implemented, as a sign or Heaviside function. Note that these are distributions, also called generalized functions⁶, and this allows us to base our solution on the representation of distributions as the *weak limit of sequences of locally integrable functions*. This has the advantage that we can select suitable locally integrable functions admitting more convenient polynomial approximations that are amenable to homomorphic encryption.

Performing the weak limit is on the other hand problematic in the homomorphic encryption setting and in general when using (polynomial) approximations, which are typically defined only over restricted intervals. We will solve this problem by selecting a class of locally integrable functions that have specific and suitable characteristics enabling us to calculate such limit in a sufficiently accurate way by mapping the values calculated over the restricted interval(s) to values at points outside such interval(s). A key-point will be keeping the number of consecutive operations sufficiently small (thus also keeping the necessary polynomials to be of a low degree).

After this general introduction to our solution, we now pass to its concrete illustration, in three key points hereby described in 1), 2) and 3).

1) Choice of the sequence of locally integrable functions.

As we mentioned, it is well-known that H can be obtained as the weak limit of several sequences of locally integrable functions, but in order to effectively perform the weak limit in homomorphic encryption with lower polynomials, we use the sequence

$$\{\tanh(kx)\}_k \quad (5)$$

where $\tanh(x)$ is the hyperbolic tangent, such that the weak limit becomes

$$H(x) = \lim_{k \rightarrow \infty} \frac{1}{2}(1 + \tanh(kx)). \quad (6)$$

The sequence of hyperbolic tangent functions will be crucial to allow us to effectively compute the weak limit in homomorphic encryption, as we will now explain.

Indeed, in homomorphic encryption we will have to polynomially approximate the functions $\tanh(kx)$. The approximation is necessarily only valid (that is, accurate) over a restricted interval. In fact for performance reasons in our case, the interval will be $[0, 0.25)$ in order to use the lowest possible order in polynomial approximations of the functions, as we will explain in point 3). However, computing the weak limit means calculating the function over large intervals defined by $z = kx, k \gg 1$. We will manage to do so precisely because of the so-called *bisection property* of $\tanh(z)$, which is why such sequence is crucial for us to be able map values calculated over

⁵We use the half-maximum convention.

⁶For an introduction see [42].

the restricted interval to values at points outside such interval and obtain the weak limit rather efficiently (only low order polynomials will be necessary).

2) Definition and calculation of the weak limit.

The weak limit in equation (6) requires mapping the (approximate) calculated value of $\tanh(z)$ for $|z| \in [0, 0.25]$ to much larger z , effectively $z = kx, k \gg 1$ (as $k \rightarrow \infty$). In order to do so, as mentioned, we employ the *bisection property*:

$$\tanh(2z) = \frac{2 \tanh(z)}{1 + \tanh^2(z)}. \quad (7)$$

After r applications of this formula, the hyperbolic tangent initially calculated at $z = x$ is now calculated at $z = 2^r x$, hence for $k = 2^r$ the limit $k \gg 1$ then corresponds to $r \gg 1$. We will discuss later on what values of r are achievable and/or efficient in practice, and what effect this has on the accuracy of the final comparison result.

The issue of divisions. Note that equation (7) involves calculating a division, which is not possible in the present homomorphic encryption schemes. We solve this issue by using specific polynomial approximations for the function $\frac{1}{x}$, where x can then be generalised to functions of our ciphertexts, once we understand how to approximate the reciprocal function for a variable x . In the case of (7), as for all z , $0 \leq \tanh^2(z) \leq 1$, we must polynomially approximate the function $\frac{1}{1+x}$ for $x \in [0, 1]$ with $x = \tanh^2(z)$. We obtain the approximation of such function by shifting $x \rightarrow x + 1$ in the approximation of the reciprocal function $\frac{1}{x}$ for $x \in [1, 2]$.

Obviously the higher the degree of the polynomial, the more accurate is the approximation, but, as said, we have to consider small degree polynomials because of the limitations of the levelled homomorphic cryptosystems we have to deal with in practice. For illustration, let us consider the two lowest minimax approximations (which can be found for example using the Matlab minimax algorithm):

$$\frac{1}{x} \approx 2.871320 - 3.029870x + 1.392785x^2 - 0.235498x^3 \quad (8)$$

$$\frac{1}{x} \approx 1.4571 - 0.5x \quad (9)$$

both for $x \in [1, 2]$. The first polynomial provides a better approximation (accuracy⁷ $\mu = 9.62$ bits, compared to 4.5 bits of the second one), but its degree is unfortunately bigger making it a less convenient candidate for homomorphic cryptosystems.

Coming back to equation (7), we thus apply $x \rightarrow x+1$ in the approximation (9) of the reciprocal function $\frac{1}{x}$ for $x \in [1, 2]$, which leads us to

$$\frac{1}{1+x} \approx 0.9571 - 0.5x \quad x \in [0, 1]. \quad (10)$$

Hence we can write the approximate bisection formula as

$$\tanh(2z) \approx \tanh(z)(1.9142 - \tanh(z)^2). \quad (11)$$

3) Polynomial approximation of the locally integrable functions.

We finally address the polynomial approximation of $\tanh(z)$ itself. Our choice for the explicit polynomial must also be guided by the fact that we are constrained in practice by the maximum number of consecutive operations that the levelled homomorphic system we are limited to can sustain before the need to bootstrap. Luckily, $\tanh(z) \sim z$ for $|z| \in [0, 0.25]$ with already quite good accuracy (≥ 7.6 bits).

In practical applications with concrete datasets, this implies that datapoints must be preprocessed and in particular *normalised* such that the values we want to compare fall within the interval $[-0.12, 0.12]$ in order to apply the algorithms with the above described approximation. **To conclude:**

Comparison Operation. For $x_1, x_2 \in [-0.12, 0.12]$, a comparison operator

$$\text{Comp}(x_1, x_2) \in \{0, \pm 1\}$$

can be implemented under homomorphic encryption by approximating the Heaviside function for $|x| < 0.25$ with

$$\lim_{r \rightarrow \infty} \frac{1}{2}(1 + \tanh(2^r x))$$

through iteratively replacing $\tanh(2x)$ by $x(1.9142 - x^2)$.

We detail the pseudocode of the comparison operation in Algorithm 1.

B. Implementing Selection/Jump Operations in Homomorphic Cryptosystems

As we have mentioned above, the selection/jump operation is *particularly difficult in an homomorphic setting* and this is a crucial realisation of our analysis. Indeed, public-key cryptography requires *ciphertext indistinguishability*, which is evidently in tension with the necessity to select a path (one or more ciphertexts) at run time, that is, before the decryption, which is supposed to occur only at the end.

We propose, as best operational solution to this issue, an “*implicit selection*” by *weighting*. This is in fact not an actual selection so that it fully respects semantic security. The idea is not to truly select, but to map the two subsets (the one of elements we “want-to-select”, and the one of elements “not-to-select”) into two different subspaces, choosing those spaces in a way that this map will keep them separate in the subsequent parts of any algorithm and will allow to recover at the end the “want-to-select” part. This is achieved by collapsing all elements of the “not-to-select” subset into the zero element of the ciphertext space, while the elements that we want to select will be preserved without change (that is, they will be mapped in themselves). We recall that in the case of homomorphic cryptosystems defined on polynomial rings the zero element is the zero polynomial.

The mapping procedure consists in re-scaling the compared data ciphertexts x_1, x_2 with suitable weights that depend on the result of the comparison. There are different ways to implement such “selection” weights, differing in what comparison operation one wants to implement ($>$, $<$, \dots) and what are the constraints on the number of consecutive operations.

⁷The accuracy μ is related to the error ϵ as $\mu = -\log_2 \epsilon$.

We now present implementations of comparison and a series of *selection* operations, each of which can be realized as an algorithm in homomorphic encryption:

Selection Operations.

$$\text{Comp} : C \times C \rightarrow \{0, \pm 1\}, (x_1, x_2) \rightarrow w_{12} \quad (12)$$

$$\text{Select}_{> \frac{1}{2}} : C \times C \rightarrow C \times C, (x_1, x_2) \rightarrow (s_{12}x_1, s_{21}x_2) \quad (13)$$

$$\text{Select}_{< \frac{1}{2}} : C \times C \rightarrow C \times C, (x_1, x_2) \rightarrow (s_{21}x_1, s_{12}x_2) \quad (14)$$

$$\text{Select}_{=} : C \times C \rightarrow C \times C, (x_1, x_2) \rightarrow (\bar{s}x_1, \bar{s}x_2) \quad (15)$$

$$\text{Select}_{>} : C \times C \rightarrow C \times C, (x_1, x_2) \rightarrow (\tilde{s}_{12}x_1, \tilde{s}_{21}x_2) \quad (16)$$

$$\text{Select}_{<} : C \times C \rightarrow C \times C, (x_1, x_2) \rightarrow (\tilde{s}_{21}x_1, \tilde{s}_{12}x_2) \quad (17)$$

with

$$s_{ij} = \frac{1 + w_{ij}}{2}, \quad \bar{s} = 1 + w_{12}w_{21},$$

$$\tilde{s}_{ij} = w_{ij} \frac{1 + w_{ij}}{2} \quad w_{ij} = -w_{ji}.$$

Note that, although $w_{21} = -w_{12}$, it is more convenient in the homomorphic cryptosystem scenario to calculate w_{12} and w_{21} independently, so that they have the same (lower) noise content, rather than w_{21} having a higher one due to being the negation of w_{12} . This will improve accuracy and precision of the algorithms allowing more operations on the ciphertexts, but at the expense of time efficiency.

The $\text{Select}_{> \frac{1}{2}}$ and the $\text{Select}_{> \frac{1}{2}}$ algorithms differ in how they map the case $x_1 = x_2$: the former map $x_{1,2} \rightarrow 0.5x_{1,2}$, the latter map $x_{1,2} \rightarrow 0$. Note that although the former algorithms do not implement exactly the $>$ and $<$ relations, they are convenient because they use less operations, and for some practical applications their treatment of the case $x_1 = x_2$ is not very problematic.

Finally, in Algorithm 2 we provide a detailed description using $\text{Select}_{> \frac{1}{2}}$ as an example (the algorithms for $\text{Select}_{> \frac{1}{2}}$ and the $\text{Select}_{> \frac{1}{2}}$ can be easily derived therefrom).

We end this section with some comments on the specific features of the mechanism we have proposed to implement the selection operations. First of all, we stress the main difference with an actual selection: while the latter operating on a certain set of elements returns in general a subset of it (typically, but of course not always, with fewer elements), our proposed mechanisms projects the unwanted elements onto the zero element, while preserving the elements one wants to select.

However, both the w_{ij} 's and the elements will be encrypted in the homomorphic case, thus we will not be able to discern which elements have been mapped to the zero element and which have been preserved (selected). Therefore, one will have to carry over all elements until the moment of decryption. This clearly has implications for the efficiency of practical applications with large datasets. We explore some of the consequences of this in Section VI-A.

We also observe that the final obtained values $s_{ij}, \bar{s}_{ij}, \tilde{s}_{ij}$ will never be exactly 1 or 0, but will tend asymptotically to

Algorithm 1

Input: Integer r and encrypted $z_c = x_{c1} - x_{c2}$, where x_{c1}, x_{c2} are encryptions of $x_1, x_2 \in [-0.12, 0.12]$ encoded using fractional encoder

Constants coefficient list $b_{\text{list}} = [-1.9142, 1.0, 0.5]$

Output: Binary values $\{0, 1\}$ with accuracy of about 3.65 bits

Algorithm

for $b \in$ coefficient list b_{list} **do do**

$b_e \leftarrow \text{Enc}_f(b)$

end for

for $i = 0$ to r **do do**

Compute: $y_c \leftarrow z_c * b_e$

Add plain: $u_c \leftarrow -1.9142e + y_c$

Multiply: $t_c \leftarrow z_c * u_c$

end for

if $r \% 2 == 1$ **then**

Negate: $z_c \leftarrow -t_c$

else

Assign: $z_c \leftarrow t_c$

end if

return $z_c = w_{12}$

Algorithm 2

Input: Integer r and x_{c1}, x_{c2} , which are encryptions of $x_1, x_2 \in [-0.12, 0.12]$ encoded using fractional encoder

Constants coefficient list $b_{\text{list}} = [-1.9142, 1.0, 0.5]$

Output: Binary values $\{0, 1\}$ with accuracy of about 3.65 bits

Algorithm

$s_{ij} = \text{Comp}(x_{ci}, x_{cj}; b_{\text{list}})$ for $ij = 12$ and 21

Add plain: $s_{ij} \leftarrow s_{ij} + 1.0e$

Multiply: $s_{ij} \leftarrow s_{ij} * 0.5e$ **return** $(s_{12}x_{c1}, s_{21}x_{c2})$

those values. An important figure of merit for the functions we have defined is the maximum number of consecutive operations they require, because the efficiency required in practical applications forces us to avoid bootstrapping, and thus allows only a limited number of consecutive operations. We will discuss this in detail in Section V-C2.

V. TESTS AND RESULTS

A. Methodology

The general results presented in this work are *agnostic* for what concerns the choice of (fully) homomorphic cryptosystem. Nevertheless, to concretely implement our models and algorithms, we have chosen to adopt the scheme of Fan and Vercauteren (FV) [25] for a series of reasons.

a) **Efficiency:** the FV scheme is an efficient implementation of the scheme in [18], one of the most remarkable second generation homomorphic systems.

b) **Comprehension of the operating range for the cryptosystem parameters:** determining the correct operating range of parameters for the various homomorphic cryptosystems is one of the active topics of research and it is unclear in all schemes. Other cryptosystems beside the Fan-Vercauteren one have been studied under this point of view, but their

good parameter ranges are much less clear than the already incomplete one in Fan-Vercauteren's, as one can for example see mentioned and discussed in [38], see also [19] when speaking of the popular scheme of [17]. The FV scheme has been subject to a few more studies and experiments, as for example can be seen in the documentation of libraries such as SEAL [52], and [13].

c) *State of software libraries*: this is the point where the FV scheme is particularly valuable, with examples such as [5], [26], [52]. In particular, SEAL [52] is evolving towards more explicit software engineering standards. We have been using its version 2.1, as latest updates have implemented modifications in the FV homomorphic scheme to improve the speed of calculations, but making it less simple to explore suitable ranges of parameters, see [52].

One important remark is that all libraries we know of do not actually implement the *fully* homomorphic cryptosystem, because they do not implement the bootstrapping, thus reducing the cryptosystem to only its somewhat homomorphic version. This will effectively limit the maximum number of consecutive operations we can evaluate.

As our work is agnostic concerning homomorphic schemes, knowing the details of the FV one is not essential. It is however relevant to have a picture of the scheme's parameters, as they affect, for instance, the size of encodable data, the size of evaluable circuits, and so on. They are:

- the plaintext modulus t , for the plaintext space $R_t \equiv \frac{\mathbb{Z}_t[x]}{f(x)}$
- the ciphertext modulus $q \gg t$, for the ciphertext space $R_q \equiv \frac{\mathbb{Z}_q[x]}{f(x)}$,
- the degree $d = 2n$ of the monic irreducible polynomial modulus $f(x) = x^d + 1$ (even degree and specific form chosen by SEAL for efficiency reasons).

We will also encode data in plaintexts using the so-called *fractional encoding* of [23], by expanding our finite precision floats u in a basis b as $u = \sum_{i=0}^u u_i b^i + \sum_{j=1}^s u_j b^{-j}$, with u_i, u_j then mapped to plaintext polynomial coefficients. This encoding depends on three parameters:

- the basis b ,
- the number of polynomial coefficients reserved for the fractional part $n_f = \max$ allowed s ,
- the number of polynomial coefficients reserved for the integer part $n_i = \max$ allowed u .

B. Datasets

As we have discussed in the introduction, only integers and finite precision floats (that is, rational numbers) are representable in the existing homomorphic cryptosystems. The datasets we have been using in our analysis are random datasets obtained from a uniform distribution of float values, and normalised according to the specifics we will illustrate in Sections IV-A and IV-B.

C. Empirical Study

We now turn to the empirical study of the algorithms leading to the functions in equations (12 - 17), which all depend on the single parameter r , the number of iterations to compute in the weak limit approximation derived from equation (7). We want to determine for which values of r and range of data

arguments x_1, x_2 the algorithm is sufficiently accurate and this will involve studying the algorithms both in their unencrypted and encrypted form.

The tests are run over different datasets, specified in the respective subsections. The evaluation of the algorithm performances are based on the *Mean Absolute Error* (MAE):

$$\text{MAE}(X, Y) = \sum_{y \in Y} \frac{|x - y|}{|Y|} \quad (18)$$

where X is the set of expected values x , Y the set of obtained ones (from the algorithms) and $|Y| = |X|$ denotes its cardinality.

1) *Evaluation of algorithm parameters - unencrypted form of the algorithm*: We first study the algorithms in unencrypted form to establish the dependence of the results on r . We have considered a set X_t of samples x in the interval $|x| \in [0, 0.25)$ where the algorithms (12 - 17) can operate. For each sample we have run the algorithm several times, for an increasing number of iterations r , starting from $r = 1$. We have then evaluated the accuracy of the algorithm in the unencrypted form by calculating the MAE.

We present in Figure 1 a series of illustrative plots. We have chosen to report here plots concerning $\text{Select}_{> \frac{1}{2}}(x, 0)$ with little loss of generality concerning the illustrative purpose, as those for $\text{Select}_{>}(x, 0)$ are quite similar, and as the algorithms for the $<$ operations are the same up to an intermediate sign. The plots in blue (rows one and three) show the returned results from $\text{Select}_{> \frac{1}{2}}(x, 0)$ against the number of iterations r . The plots in red (row two and four) show the value of the simple error defined as

$$\text{Simple_Error}(x, r) = H(x) - \text{Select}_{> \frac{1}{2}}(x, 0) \quad (19)$$

where $H(x)$ is the Heaviside function. We have chosen to plot the results for some of the values x we have considered. In particular we plot for points with values x going from -0.20 to 0.20 in steps of 0.05 in order to provide a consistent coverage of the working interval of the algorithms.

We also present in Table I the values of the MAE, equation 18, for those set of values for x and for the tested number of iterations r in the cases of $\text{Select}_{> \frac{1}{2}}(x, 0)$ and $\text{Select}_{>}(x, 0)$. The results in the table show that, in the former, for $r = 4$ the MAE has dropped at around 6%, and for $r = 5$ around 3%, while in the latter the values are slightly bigger. The low degree of the approximating polynomial from equation (11), and thus low number of necessary consecutive multiplications, make this an interesting result for applications in homomorphic encryption.

2) *Selection of algorithm and homomorphic scheme parameters - encrypted form of the algorithm*: We move now to a series of tests with a carefully chosen artificial dataset to establish in the best r and FV cryptosystem parameters, where *best* means leading to the smallest errors over the maximum possible data-value interval.

The value of r determines the number of consecutive operations the algorithms *must* sustain, while the parameters

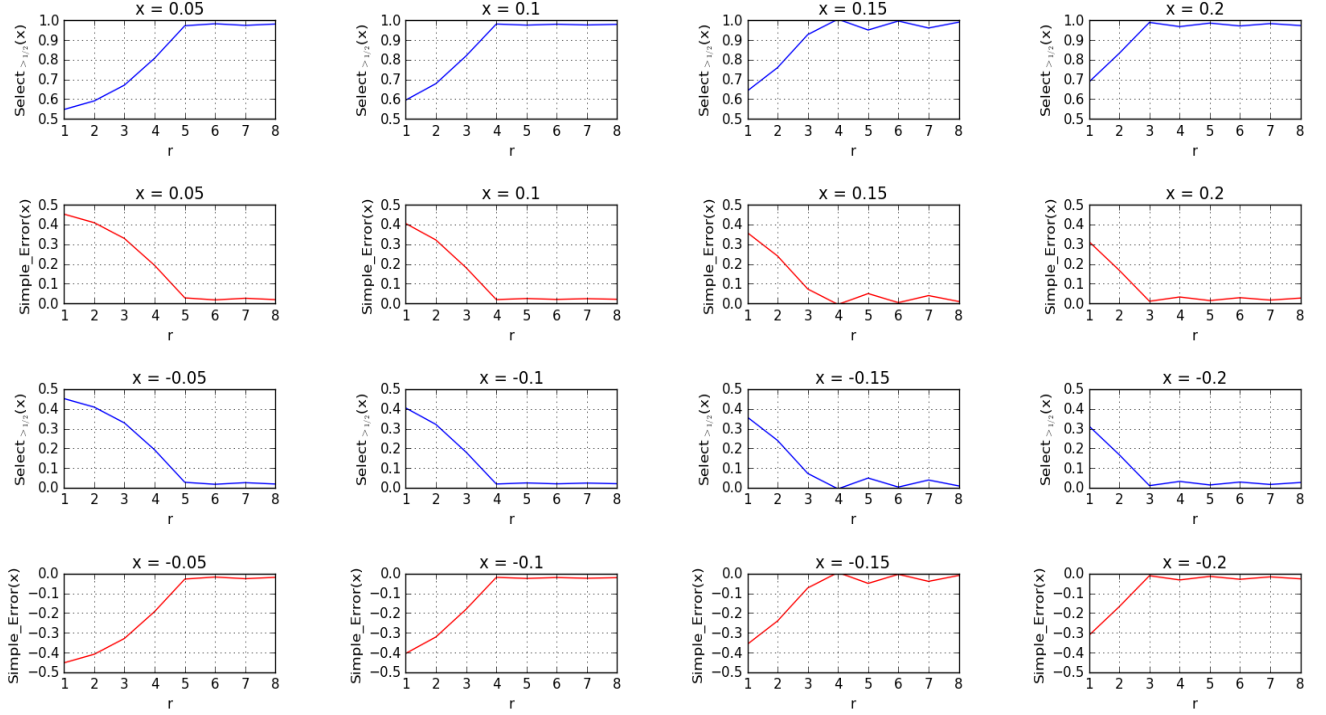


Fig. 1: Plots in rows one and three (blue color) show the value of the returned results from $\text{Select}_{> \frac{1}{2}}(z, 0)$ against the number of iterations r . The plots in row two and four (red color) show the value of the simple error defined in Equation 19.

	$r = 1$	$r = 2$	$r = 3$	$r = 4$	$r = 5$	$r = 6$	$r = 7$	$r = 8$
MAE(r) for $\text{Select}_{> \frac{1}{2}}$	0.38	0.28	0.15	0.062	0.027	0.017	0.026	0.019
MAE(r) for $\text{Select}_{>}$	0.47	0.39	0.22	0.10	0.056	0.034	0.051	0.036

TABLE I: Values for the mean absolute error defined in Equation 18 for a number of iterations values r calculated over a set of points $x_t = 0.05 * s$, $-4 \leq s \leq 4$.

of the homomorphic scheme determine the number of consecutive operations the *encrypted* algorithm *can* sustain without bootstrapping. In the case of algorithms (13 - 17), the number of (noise-dominating) consecutive multiplications we need to be able to perform to run r iterations is

$$2r + 1_p \quad \text{for the } \text{Select}_{> / < \frac{1}{2}} \text{ algorithms} \quad (20)$$

$$2r + 1 + 1_p \quad \text{for the } \text{Select}_{> / \leq} \text{ algorithms} \quad (21)$$

where 1_p is a multiplication with a plaintext coefficient⁸, and $2r$ or $2r + 1$ are ciphertext multiplications. The total count of operations is of course higher, when including additions and relinearizations, but as they generate less noise, we will neglect them. Moreover, note that the ciphertext multiplications involve multiplying recursively the *same* ciphertext, which means that successive multiplications are more costly for the noise growth, as they involved already noise-grown ciphertexts.

A number of consecutive operation such as, for instance, $8+1$ (for $r = 4$) may seem not huge, but it must be put in

⁸We distinguish mixed plaintext/ciphertext multiplications because the noise level estimates are different than pure ciphertext multiplications in our implementation based on SEAL, see Table 3 in [19].

relation with the cryptosystem parameters necessary to accommodate it. As said before, the analysis concerning the choice of parameters is still an active field of research, and there are in fact different partial results in the literature. For example [20] estimates the cryptosystem parameters for a scheme like the FV one, finding that already to perform 10 multiplications (of different, and thus with minimal noise, ciphertexts) requires a polynomial modulus degree of $d = 8192$ and a plaintext modulus of at least 2^{243} for a fractional encoding in base $b = 3$, which in turns implies a value for the ciphertext modulus of about 2^{226} to have 123 bits security as estimated in [19]. However, the work [6] claims that much higher values are actually necessary to be able to perform 4 subsequent multiplications, already in the case of a much simpler integer encoding ($t = 131702$, $q > 2^{159}$, $d = 81920$). Finally, the recent paper [15] claimed necessary values of the form $t \gtrsim 2^{107}$, $d \leq 368$ when dealing with the case of a graph (representing an instance of Ivakhnenko's group method of data handling), whose evaluation along a path from input to output comported ≈ 6 consecutive multiplications (plus a similar quantity of additions).

In summary, two things appear from the literature:

- the parameter choice bounds are coarsely estimated for similar experiments,
- the number of consecutive multiplications implemented in existing literature is very low, thus our $\sim 8 + 1$ one appears to be the highest ever tried.

We now present the results for the algorithms, $\text{Select}_{>/< \frac{1}{2}}$ and $\text{Select}_{>/<}$ (given the number of operations required by these, the results also apply to the case $\text{Select}_{=}$). We have been testing with parameters ranging over a number of possible values, in particular

$$\begin{aligned}
 d &\in \{8192, 16384, 32768\} \\
 q &\in \{2^{116} - 2^{18} + 1, 2^{226} - 2^{26} + 1, \\
 &\quad 2^{435} - 2^{33} + 1, 2^{829} - 2^{54} - 2^{53} - 2^{52} + 1\} \\
 t &\in \{4096, 16384, 65536\} \\
 b &\in \{3, 5, 7, 9\} \\
 n_f &\in \{6, 8, 10, 24, 32\} \\
 n_i &\in \{8, 16, 32, 64\}
 \end{aligned} \tag{22}$$

where the parameters and their notation have been defined in Section V-A. The time of key generation and storage overhead following from these choices of parameters are known for the FV scheme and the SEAL implementation library, which we have used in our tests, and we refer the reader to the literature, see the original articles [25], [52]. We also recall, and stress, that our algorithms are agnostic for what concerns the choice of homomorphic scheme, and thus any performing scheme may be used in practical applications, so that computation and storage overheads can be tuned as desired.

The range of values for q, t, d were chosen by taking advantage of the sets of values indicated by the SEAL team in their testings of the library versions 2.1 and 2.2. Instead, SEAL version 2.3.0 uses a modification of the FV scheme to increase time efficiency, but which allow somewhat less flexibility and “ease” in the choice of parameters. In particular, the available values of the parameter q in version 2.3.0 have proven in our case to yield sub-optimal results.

We have run the algorithm $\text{Comp}(z, 0)$ over a small dataset

$$z \in \{-0.20, -0.15, -0.10, -0.05, 0, 0.05, 0.10, 0.15, 0.20\}$$

capable however to cover sufficiently uniformly the allowed instance space $|z| \in [0, 0.25)$ from Sections IV-A and IV-B.

We list in Table II the best results for each tested value of r , where best indicates smallest MAE for the selection weights $s_{ij}, \bar{s}_{ij}, \tilde{s}_{ij}$ as defined in Equation (18).

From our analysis, the rationale behind the effectiveness of encryption scheme parameters emerges as follows. First of all we need small t and large q since $\frac{q}{t}$ mostly determines the maximum noise bound, see [19]. Secondly, we need to keep the number of coefficients reserved to the fractional part in encoding (n_f) as small as possible because during multiplication the number of coefficients occupied by the fractional part will increase rapidly. The number of coefficients reserved to the integer part (n_i) is of less concern, because all the normalised test data instances x are smaller than one.

The basis b used for the fractional encoding, see Section V-A also played an important role. One would like to have as small a basis as possible, to avoid the “wrapping up” of the

modulo t during computations. However, smaller basis also means that more coefficients of the plaintext polynomial will be non-zero, and so since the number of coefficients of the fractional part increase with multiplications, they can more easily cross over to the coefficients reserved for the integer part, and ruin decryption.

Finally, the degree d of the polynomial modulus is relevant because of two different reasons: on the one hand the experiment should take into account security bounds, which depend on d since long polynomial are more difficult to attack; on the second hand having a big number of coefficients also helps in avoiding that those reserved to the fractional part and those to the integer part cross over and mix up rapidly.

3) **Full tests - encrypted form of the algorithm:** Having estimated as discussed the best algorithm and scheme parameters, we have finally run full-fledged tests over randomised datasets to assess the accuracy of the algorithms.

We have studied the algorithms $\text{Select}_{> \frac{1}{2}}$, $\text{Select}_{>}$ and $\text{Select}_{=}$, since the algorithms for the $<$ (less than) relations are essentially the same as the ones for $>$ up to an (intermediate) sign, hence the test results apply to those as well. Our datasets consisted of couples of datapoints randomly generated in the range $[-0.12, 0.12]$ (so that the difference between datapoint values would fall in the valid range $[0, 0.25)$ to apply the algorithms, see Sections IV-A and IV-B). We have used several measures of accuracy and performance for the algorithms, to be able to provide a rigorous evaluation.

The results are presented in Table III. The simplified notation $\text{MAE}(a)$ indicates the error calculated using equation (18) on the values $a = a_{\text{out}} - a_{\text{expected}}$. We have studied various tests, in particular we have considered a to be first $s_{ij}, \bar{s}_{ij}, \tilde{s}_{ij}$ and then the final full output of the algorithms (that is, $a = s_{ij}x_i$ and the analogous for $\bar{s}_{ij}, \tilde{s}_{ij}$).

The best performing algorithms are $\text{Select}_{> \frac{1}{2}}$ (and thus $\text{Select}_{< \frac{1}{2}}$ as it is the same up to an initial sign), achieving about 20% error on the selection weights and 2% on the final output. The error is dominated by the error value for datapoints that are very close to each other. In fact, we have run the same tests with datapoints with a fixed minimal distance in order to check variations depending on this, and the error rate drops rapidly in function of the inter-distance of points (already with inter-distance higher than a few percent, for instance 3%, the error rate on selection weights drop at about 12%).

The algorithm $\text{Select}_{=}$ deserves a special comment: the comparison weights w_{ij} are exactly 1 when $x_i = x_j$ and different from 1 when $x_i \neq x_j$ (and closer to 0 as the difference/distance between x_i and x_j is larger), so that if in a simple application one lets the data owner simply decrypt the comparison weights and pick the datapoints corresponding to weights equal to 1 to perform the selection part of the conditionals of interest, one would have perfect accuracy. This however cannot be done when the algorithm must be inserted in a longer pipeline of algorithms and the “selection” must be performed on the encrypted parts and carried over to further steps of the pipeline. The results relative to $\text{Select}_{=}$ that we show in Table III are therefore to be intended for this case.

We finally present in Table IV the result for the timing of

Iterations	Results
$r = 3$	Smallest parameters where result still achieved: $d = 16384, q = 2^{435} - 2^{33} + 1, t = 65536, b = 7, n_i = 8, n_f = 8$
$r = 4$	Not accomplished correctly (error less than 1 at least) by any parameter value in 22. "Best results" for $d = 16384, q = 2^{435} - 2^{33} + 1, t = 65536, b = 7, n_i = 8, n_f = 8$

TABLE II: Values for the mean absolute error defined in Equation 18 for a number of iterations values r calculated over a set of points $x_t = 0.05 * s, -4 \leq s \leq 4$.

$d = 16384, q = 2^{435} - 2^{33} + 1, t = 65536, b = 7, n_i = 8, n_f = 8$						
Iterations	Select $_{>\frac{1}{2}}$		Select $_{>}$		Select $_{=}$	
	MAE(s_{ij})	MAE($s_{ij}x$)	MAE(\tilde{s}_{ij})	MAE($\tilde{s}_{ij}x$)	MAE(\tilde{s}_{ij})	MAE($\tilde{s}_{ij}x$)
$r = 3$	0.26	0.021	0.41	0.023	0.52	0.057
$r = 4$	1.7	0.28	4.9	0.55	2.6	0.30

TABLE III: Errors for the comparison and selection/jump algorithms defined in (12), (13), (15), (16). We have tested the algorithms on randomly generated datasets with batches of 60 couples of datapoints to be compared and present here the average results for $r = 3$, while for $r = 4$ we present the result for the best batch (since anyway the case $r = 4$ is affected by error of decryption due to too many consecutive operations performed, see Section V).

Average timing per instance in seconds			
$d = 16384, q = 2^{435} - 2^{33} + 1, t = 65536, b = 7, n_i = 8, n_f = 8$			
Iterations	Select $_{>\frac{1}{2}}$	Select $_{>}$	Select $_{=}$
$r = 3$	17.4 s	21.5 s	21.1 s
$r = 4$	30.5 s	31.5 s	31.2 s

TABLE IV: Values for the timing for runs of the selection algorithms in seconds per instance.

the selection/jump algorithm. Our work has not focused on achieving the best performance, as it has been more centred on the proof-of-concept and the practical implementation of the algorithms, as well as to the discussion of the novel issues concerning homomorphic encryption and applications such as machine learning (see Section VI-A). We have however measured the timings when running our tests, and report in the table the average timing per (x_1, x_2) instance for the algorithms (13), (15), (16). A direct comparison with the results reported in the literature is however not straightforward, because:

- there are very few works implementing similarly complex algorithms in a homomorphic cryptosystem,
- often the experiments in the literature have been performed on powerful computer clusters, see e.g. [6],
- only few among the works with complex algorithms report full algorithm timings⁹.

D. Improvements and comments

We have presented here above a series of algorithms to evaluate comparisons and conditionals in homomorphic encryption settings. The algorithms have been explicitly tested in a concrete implementation of the Fan-Vercauteren encryption scheme in a levelled form. The limitation on the total number of consecutive operations has the strongest influence on the accuracy of the algorithms.

Such limitations, and hence inaccuracies, would simply be absent for an implementation in a fully homomorphic scheme,

⁹The others, e.g. [15], report "time per operation" such as addition, multiplication or encryption. However, also other routines such as relinearizations are present and finding the overall algorithm timing is not straightforward.

or, possibly, using schemes that although limited can tolerate a larger number of consecutive operations (we estimated 9 multiplications would already guarantee accuracy at percent or sub percent level).

It would be also interesting to assess what effects new plain- and ciphertext encodings such as [13] would have, possibly in alleviating some of the accuracy loss due to crossing of the integer and fractional parts of the standard encoder, see Section V-C3.

VI. APPLICATIONS

A. Machine learning and specific issues

As mentioned in the introduction, part of the present interest in homomorphic encryption stems from the potential to permit privacy preservation to coexist with the nowadays ubiquitous machine learning/data mining/predictive analytics¹⁰ without incurring in the limits of the privacy/data usefulness budget of other approaches [28].

In the literature of the past few years we can find a limited number of implementations of machine learning algorithms¹¹ preserving privacy combined with homomorphic cryptography techniques, see e.g. [6], [13], [14], [32], [35].

The main problem is that very often machine learning techniques are not amenable to homomorphic encryption due to various limitations and it is therefore interesting to re-think the actual machine learning algorithms.

¹⁰Again, for brevity of expression in the following we will use "machine learning" as an umbrella term for all these different but related approaches.

¹¹Typically for prediction only, that is having the training part all in unencrypted form.

1) *Relevant aspects of machine learning (ML):*

Developing machine learning/data mining/predictive analysis systems typically involves a series of different steps¹²: training, validation, testing, prediction. The core issue is coined as solving a “learning problem”, which comes down to finding within a certain solution space (essentially delimited by the inference bias) a function or generalisation thereof that maps input data to a correctly inferred output (be it classification or regression). To this end, the algorithm uses the input data to assess the relevance of different hypothesis in the solution space, building them against the available training data, and validating and testing them against independent pieces of data. The tested algorithm can then be used with other data for prediction.

2) *Problematic points of ML in homomorphic settings:*

We will now elaborate on certain ingredients of the machine learning inference process that clash in a particularly relevant way with fundamental constraints of homomorphic encryption. In a large class of machine learning algorithms two element are paramount: the *stopping criterion* and *heuristics*.

- **Stopping criterion.** Typically machine learning algorithms terminate when a stopping criterion is met, generally when an extremal condition is reached. This means that the algorithm must be able to evaluate a condition (the criterion) and select one of the options (essentially, continue or stop) while running in its encrypted form. In Section IV-B we explained that the selection step conflicts with the fundamental requirement of semantic security in (homomorphic) encryption and we proposed a “selection by weighting” algorithm that allows mapping the selected-for and unselected-for subsets to specific subsets (in particular the zero element for the unselected subset) that at decryption will provide the desired result. The bigger problem in this case is that the “selection by weighting” does not really signal that a selection has been made, but all mapped results (both “selected” and “rejected” ones) are still encrypted and carried over. There is no way at run time to determine that the stopping criterion has thus been met and the procedure must be stopped, until decryption occurs. Unfortunately, and importantly, not stopping the training of an algorithm precisely at the stopping point does entail overfitting and thus suboptimal learning models.
- **Heuristics.** In order to efficiently explore the instance (data) and problem space, and make useful inference, several machine learning algorithms operate heuristic choices at run time. Again, the clash between the need to make selections and the requirement of semantic security of the (homomorphic) encryption pops up. Differently from the case of the stopping criterion, here our “selection by weighting” would not create loss of accuracy in the algorithms, but, of course, in the case of large datasets it would entail carrying over the full dataset all along, hence affecting the efficiency of the algorithm, and in

certain cases its whole inference capabilities, as we will discuss at the end of this section.

Also another issue can arise: the comparison weights w_{ij} are not exactly 1 or 0 but some other numbers (float) close to that, because of the limitations in the number of consecutive operations of the levelled homomorphic system one has to use in practice, which limits the value of the algorithm parameter r and thus its accuracy. This effectively transforms a machine learning algorithm into a weighted version of itself. In some cases this does not significantly affect the accuracy, sensitivity and precision of the algorithms, but in other cases it does, also in an adverse way. The studies in this respect are scarce in the literature, see e.g. [1] for what concerns clustering.

The two issues here presented are quite fundamental and could seriously complicate, if not make impossible, the implementation of privacy-preserving machine learning using purely homomorphic encryption. The stopping criterion issue, in particular, implies that the training of correctly performing algorithms (that is, not overfitted ones) does not seem achievable without decryption at run time, which goes against the aim itself of homomorphic encryption. While this drawback affects only the training of models, private data are also used in training (even more than in the prediction runs after training) and should therefore be protected as well.

The heuristics issue instead seems to represent a secondary problem, only affecting performance and thus possibly solved by, for example, hardware evolution. However, that is not the case. To be able to make actual inferences several machine learning algorithms do need to operate with heuristics on the data/problem space. If that is not possible, the algorithms cannot proceed with meaningful inferences. Again, this appears to be a relevant obstacle on the road to make machine learning privacy-friendly by using homomorphic encryption.

B. *Applications to algorithms different than machine learning*

Algorithms different from machine learning or similar predictive analytic techniques that do not need to make inferences and avoid overfitting as discussed in the previous section, are not in such a relevant clash on general grounds with homomorphic encryption.

This means that operations such as pure database searches, for instance, even including comparisons, conditionals and selections could be effectively performed taking advantage of the techniques and algorithms we have developed in this work and their future improvements.

VII. CONCLUSIONS

Homomorphic encryption provides in theory an elegant solution to the problem of privacy preservation in data-based applications, such as those provided and/or facilitated by machine learning techniques, but several limitations hamper its implementation. In this work we have identified, on the basis of the structured program theorem, the set of minimal operations that guarantee the computation of any computable function or algorithm. We have then focused on those that are still lacking in homomorphic encryption, namely comparisons and

¹²Not always: for example instance-based methods, such as k -nearest neighbours, do not need an actual training, validation and testing phase.

conditional selections. We have discovered rather fundamental clashes between the necessity to implement those operations and the basic requirements of (homomorphic) encryption. We have also proposed practical implementations for those operations or their closest possible forms in homomorphic encryption. The limitation on the total number of consecutive operations, due to the use of levelled homomorphic encryption schemes without using bootstrapping (a practical limitation we have to face), has had the strongest influence on the accuracy of our algorithms. Percent accuracy (and better) can be obtained however for datapoints which are sufficiently interspaced. Moreover, such limitations, and hence inaccuracy, would not occur in a fully homomorphic scheme, or, possibly, using schemes and/or encodings that can tolerate a larger number of consecutive operations even if only somewhat homomorphic (we estimate from 9 multiplications onward).

We have also analysed the specific situation arising in machine learning/predictive analytic applications. We have pointed out at least two main sources of tension with the use of homomorphic encryption to fully guarantee privacy preservation in machine learning due to the newly found above-mentioned issues. These two sources of tensions are the stopping criterion and heuristics. They are present and paramount in most machine learning algorithms, and clash with (homomorphic) encryption in that they require performing selection/jump operations at run time, which in its turn clashes with semantic security, as we have studied in this work.

Two options are open under this respect. On the one hand it might be possible to find new classes and families of learning algorithms that operate without “choices at run time”. On the other, we could reconsider the use of homomorphic encryption. Maybe some other technology, such as for example functional encryption (e.g. [29]) may be capable to avoid the need for high level of communication and intermediate decryption of other techniques (such as secure multi-party computation ones)? Functional encryption could allow to limit the computations to some agreed level, while preserving the rest of the privacy of the data or algorithm.

Note however that the class of algorithms which do not need to make inferences and where overfitting can be avoided, are not in conflict with the general grounds of homomorphic encryption. We believe that further exploring the use of homomorphic encryption in algorithms for privacy preservation is paramount.

VIII. REFERENCES

- [1] M. Ackerman, S. Ben-David, S. Branzei, D. Loker, “Weighted clustering”. In Proc. 26th AAAI Conference on Artificial Intelligence, 2012.
- [2] C. Aggarwal, P. Yu, “Privacy-Preserving Data Mining: Models and Algorithms”, Kluwer Academic Publishers Boston/Dordrecht/London, 2008.
- [3] R. Agrawal, J. Kiernan, R. Stikant, Y. Xu, “Order preserving encryption for numeric data”, Proceedings of the 2004 ACM SIGMOD international conference on Management of data, 563-574.
- [4] F. Armknecht, C. Boyd, C. Carr, K. Gjøsteen, A. et al “A Guide to Fully Homomorphic Encryption”, Cryptology ePrint Archive, Report 2015/1192, 2015.
- [5] L. J. M. Aslett (2014), HomomorphicEncryption: Fully Homomorphic Encryption. R package version 0.2. URL: <http://www.louisaslett.com/HomomorphicEncryption/>
- [6] L. J. M. Aslett, P. M. Esperança, C. C. Holmes, “Encrypted statistical machine learning: new privacy preserving methods”, Technical report, Univ. of Oxford, 2015.
- [7] B. Barak, O. Goldreich, R. Impagliazzo, S. Rudich, A. Sahai, S. Vadhan, K. Yang, “On the (Im)possibility of Obfuscating Programs”, Advances in Cryptology — CRYPTO 2001: Proceedings 21st Annual International Cryptology Conference, Santa Barbara, California, USA, August 19–23, 2001.
- [8] F. Bêtul Durc, T.M. DuBuisson, D. Cash, “What Is revealed by Order-revealing encryption?”, ACMCCS ’16.
- [9] C. Bohm, G. Jacopini, “Flow Diagrams, Turing Machines and Languages with Only Two Formation Rules”. Communications of the ACM. 9 (5): 366–371 (May 1966).
- [10] A. Boldyreva, N. Chenette, Y. Lee, A.O. Neill, “Order-preserving symmetric encryption”, in: A. Joux (Ed.), Advances in Cryptology EUROCRYPT 2009, Vol. 5479 of LNCS, pp. 224–241.
- [11] A. Boldyreva, N. Chenette, A.O. Neill, “Order-preserving encryption revisited: improved security analysis and alternative solutions”, Advances in Cryptology, CRYPTO 2011, Vol. 6841 of LNCS, pp. 578-595.
- [12] D. Boneh, A. Sahai, B. Waters, “Functional encryption: definitions and challenges”, In proceedings of TCC’11, LNCS 6597, pp. 253-273. eprint.iacr.org/2010/543.pdf
- [13] C. Bonte, C. Bootland, J. W. Bos, W. Castryck, I. Iliashenko, and F. Vercauteren, “Faster Homomorphic Function Evaluation using Non-Integral Base Encoding,” In Cryptographic Hardware and Embedded Systems - CHES 2017, LNCS 10529, W. Fischer, and N. Homma (eds.), Springer-Verlag, pp. 579-600, 2017.
- [14] J. W. Bos, K. Lauter, M. Naehrig, “Private predictive analysis on encrypted medical data”, Journal of Biomedical Informatics 50, 234243, 2014.
- [15] J.W. Bos, W. Castryck, I. Iliashenko, F. Vercauteren, “Privacy-friendly forecasting for the Smart Grid using Homomorphic Encryption and the Group Method of Data Handling”, 2016, Cryptology eprint archive: 2016/1117.
- [16] R. Bost, R. Ada Popa, S. Tu, S. Goldwasser, “Machine Learning Classification over Encrypted Data, NDSS 2015, Cryptology ePrint Archive, Report 2014/331.
- [17] Z. Brakerski, C. Gentry, V. Vaikuntanathan, “(levelled) Fully homomorphic encryption without bootstrapping”, in Proceedings of the 3rd Innovations in Theoretical Computer Science Conference, ACM, pp. 309-325, 2012.
- [18] Z. Brakerski, “Fully homomorphic encryption without modulus switching from classical GapSVP”, In Advances in Cryptology CRYPTO 2012. Springer, 868–886, 2012.
- [19] H. Chen, K. Laine, P. Player, “Simple Encrypted Arithmetic Library - SEAL v2.1”, Cryptology ePrint Archive, Report 2017/224, <https://eprint.iacr.org/>.
- [20] N. Chenette, K. Lewi, S.A. Weis, D.J. Wu, “Practical order-revealing encryption with limited leakage”, 2015, Cryptology eprint archive: 2015/1125.

- [21] J.H. Cheon, J. Kim, M.S. Lee, A. Yun, "CRT-based fully homomorphic encryption over the integers", *Information Sciences*, Volume 310, 20 July 2015, Pages 149-162.
- [22] D. Demmler, T. Schneider, M. Zohner. 2015. "ABY-A Framework for Efficient Mixed-Protocol Secure Two-Party Computation". In 22nd Annual Network and Distributed System Security Symposium, NDSS 2015, San Diego, California, USA, February 8-11, 2015.
- [23] N. Dowlin, R. Gilad-Bachrach, K. Laine, K. Lauter, M. Naehrig, J. Wernsing, "Manual for Using Homomorphic Encryption for Bioinformatics", *Proceedings of the IEEE*, Volume: 105, Issue: 3, March 2017.
- [24] Z. Erkin, M. Franz, J. Guajardo, S. Katzenbeisser, I. Lagendijk, and T. Toft, "Privacy-Preserving Face Recognition". *PETS 2009*: 235-253.
- [25] J. Fan, F. Vercauteren, "Somewhat practical fully homomorphic encryption", *IACR Cryptology ePrint Archive* eprint.iacr.org/2012/144.
- [26] <https://github.com/CryptoExperts/FV-NFLlib/blob/master/FV.hpp>
- [27] C. Fontaine, F. Galand, "A survey of homomorphic encryption for nonspecialists", *EURASIP J. on Information Security archive*, Volume 2007, No. 15., January 2007.
- [28] M. Fredrikson, E. Lantz, S. Jha, S. Lin, D. Page, T. Ristenpart, "Privacy in Pharmacogenetics: An End-to-End Case Study of Personalized Warfarin Dosing", 23rd *USENIX Security Symposium*, 2014.
- [29] S. Garg, C. Gentry, S. Halevi, M. Raykova, A. Sahai, B. Waters, "Candidate Indistinguishability Obfuscation and Functional Encryption for All Circuits". *SIAM J. Comput.* 45(3): 882-929, 2016.
- [30] C. Gentry, "A fully homomorphic encryption scheme", PhD thesis, Stanford University, 2009.
- [31] C. Gentry, A. Sahai, B. Waters, "Homomorphic encryption from learning with errors: Conceptually-simpler, asymptotically-faster, attribute-based", in *Advances in Cryptology CRYPTO 2013*, Springer, pp. 7592, 2013.
- [32] R. Gilad-Bachrach, N. Dowlin, K. Laine et al, "CryptoNets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy", *Proc. of The 33rd Int. Conf. on Machine Learning*, pp. 201–210, 2016.
- [33] S. Goldwasser, S. Micali, "Probabilistic encryption & how to play mental poker keeping secret all partial information", *Proceedings 14th annual ACM symposium on Theory of computing*, p.365-377, May 05-07, 1982.
- [34] O. Goldreich, "A uniform complexity treatment of encryption and zero-knowledge," *Journal of Cryptology*, vol. 6, no. 1, pp. 21-53, 1993.
- [35] T. Graepel, K. Lauter, M. Naehrig "ML Confidential: Machine learning on encrypted data", in T. Kwon, M.-K. Lee and D. Kwon, eds, *Information Security and Cryptology (ICISC 2012)*, Vol. 7839 of *Lecture Notes in Computer Science*, Springer, pp. 1–21.
- [36] P. Grubbs, K. Sekniqi, V. Bindschaedler, M. Naveed, T. Ristenpart, "Leakage-Abuse Attacks against Order-Revealing Encryption", *IEEE 2017 Symposium on Security and Privacy* : 655-672, 2017.
- [37] S. Halevi, V. Shoup, (2014), *Helib*, <https://github.com/shaih/HELlib>.
- [38] S. Halevi, V. Shoup, "Bootstrapping for helib", *Advances in Cryptology - EUROCRYPT 2015 - 34th Annual Int. Conf. on the Theory and Appl. of Cryptographic Techniques*, Volume 9056 of *LNCS*, pp. 641–670, 2015.
- [39] S. Halevi, Y. Lindell, "Homomorphic Encryption", chapter in the book "Tutorials on the Foundations of Cryptography: Dedicated to Oed Goldreich.", pp. 219-276, 2017.
- [40] E. Hesamifard, H. Takabi, M. Ghasemi, R.N. Wright, *Privacy-preserving machine learning as a service*, *Proceedings on Privacy Enhancing Technologies* 2018 (3), 123-142.
- [41] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan. *GAZELLE: A Low Latency Framework for Secure Neural Network Inference*. *Usenix Security 2018*: 1651-1669.
- [42] R.P. Kanwal, "Generalized Functions: Theory and Technique", 2nd ed. Boston, MA: Birkhauser, 1998.
- [43] K. Lewi, D. Wu, "Order-Revealing Encryption: New Constructions, Applications, and Lower Bounds", *Proc. of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1167-1178, 2016.
- [44] J. Liu, M. Juuti, Y. Lu, N. Asokan. "Oblivious Neural Network Predictions via MiniONN Transformations". *ACM CCS 2017*: 619-631.
- [45] A. Lopez-Alt, E. Tromer, V. Vaikuntanathan "On-the-fly multiparty computation on the cloud via multikey fully homomorphic encryption", *Proc. of the 44th Symposium on Theory of Computing Conf. 2012*, pp. 1219–1234.
- [46] C. Mavroforakis, N. Chenette, A.O. Neill, G. Kollios, R. Canetti, "Modular Order-Preserving Encryption, Revisited", *Proc. of the 2015 ACM SIGMOD International Conf. on Management of Data*, pp. 763-777, 2015.
- [47] P. Mohassel, Y. Zhang. "SecureML: A System for Scalable Privacy-Preserving Machine Learning". *IEEE S&P 2017*: 19-38, 2017.
- [48] R. Popa, F. Li, N. Zeldovich, "An ideal-security protocol for order-preserving encoding", in: *Security and Privacy, 2013 IEEE Symposium on*, vol. 465, pp. 463-477 .
- [49] E. Ya. Remez, *Sur la détermination des polynômes d'approximation de degré donnée*, *Comm. Soc. Math. Kharkov* 10, 41,1934.
- [50] R. L. Rivest, L. Adleman, M. L. Dertouzos "On data banks and privacy homomorphisms", *Foundations of Secure Computation* 4(11), 169–180, 1978.
- [51] W. Rudin, *Principles of mathematical analysis*, McGraw-Hill, 1976.
- [52] SEAL library <https://www.microsoft.com/en-us/research/project/simple-encrypted-arithmetic-library/>.
- [53] N.P. Smart and F. Vercauteren. Fully homomorphic encryption with relatively small key and ciphertext sizes. *Public Key Cryptography PKC 2010, Lecture Notes in Comput. Sci.* 6056, 420–443, 2010.
- [54] D. Stehlé, R. Steinfeld, "Faster fully homomorphic encryption", in *Advances in Cryptology-ASIACRYPT 2010*, Springer, pp. 377394, 2010.
- [55] M. Togan, C. Plasca, "Comparison-Based computations over fully homomorphic encrypted data", *COMM 2014 International Conference*.
- [56] Wee H., "Functional Encryption and Its Impact on Cryptography", *SCN 2014: Security and Cryptography for Networks* pp 318-323, 2014.